

## A APPENDIX

### A.1 Notations

Table 5: Frequently used notations in this paper.

Notations	Descriptions
$\mathcal{G}^{(t)}$	Directed graph at time $t$
$n^{(t)}, m^{(t)}$	Numbers of nodes and edges at time $t$
$F$	Dimension of the node attribute
$X^{(t)}$	The node attribute matrix and $X^{(t)} \in \mathbb{R}^{n^{(t)} \times F}$
$\mathbf{x}_i^{(t)}$	The attribute vector of the $i$ -th dimension
$\mathcal{N}_{out}^{(t)}(v), \mathcal{N}_{in}^{(t)}(v)$	The out-neighbors and in-neighbors of $v$
$e_t = \{u_t, v_t\}, \Gamma$	The update event and the event set
$\mathbf{r}^{(t)}$	Residue vector at time $t$ and $\mathbf{r}^{(t)} \in \mathbb{R}^{n^{(t)} \times 1}$
$\mathbf{h}^{(t)}$	Reserve vector at time $t$ and $\mathbf{h}^{(t)} \in \mathbb{R}^{n^{(t)} \times 1}$
$H^{(t)}$	The node embedding matrix at time $t$
$M^{(t)}$	The state matrix at time $t$
$\bar{\mathcal{A}}, \bar{\mathcal{B}}, C$	The network parameters in SSM
$\alpha$	Teleport probability of random walks
$\lambda$	Threshold in to control the embedding distance

### A.2 Proof of Lemma 1

PROOF. According to the definition of 1-Norm of the matrix, we have:

$$\begin{aligned} \|\mathbf{H}^{(t+1)} - \mathbf{H}^{(t)}\|_1 &= \max_{f \in F} \sum_u |\mathbf{h}^{(t+1)}(u) - \mathbf{h}^{(t)}(u)| \\ &\leq \max_{f \in F} \sum_u |z^{(t+1)}(u) - z^{(t)}(u)| + 2\epsilon \end{aligned}$$

We can express  $Z^{(t)}$  as the result of global smoothness [81, 82] which satisfies:

$$Z^{(t)} = \left(I + cL^{(t)}\right)^{-1} X^{(t)}, \quad (6)$$

where  $c = \frac{1}{\alpha} - 1$  and  $L^{(t)} = I - P^{(t)}$  is the Laplacian matrix and denotes the topological status of graph  $\mathcal{G}^{(t)}$ . Hence we have:

$$\begin{aligned} \|\mathbf{H}^{(t+1)} - \mathbf{H}^{(t)}\|_1 &\leq \max_{f \in F} \left\| \left( \left(I + cL^{(t+1)}\right)^{-1} - \left(I + cL^{(t)}\right)^{-1} \right) \mathbf{x}^{(t)} \right\|_1 + 2n\epsilon \\ &\leq c \max_{f \in F} \left\| \left( \left(I + cL^{(t+1)}\right)^{-1} \left(L^{(t+1)} - L^{(t)}\right) \left(I + cL^{(t)}\right)^{-1} \right) \mathbf{x}^{(t)} \right\|_1 \\ &\quad + 2n\epsilon \leq c \max_{f \in F} \left\| \left( P^{(t+1)} - P^{(t)} \right) \mathbf{x}^{(t)} \right\|_1 + 2n\epsilon \\ &\leq \frac{1-\alpha}{\alpha} \|\Delta P \cdot \mathbf{x}_{max}\|_1 + 2n\epsilon, \end{aligned}$$

where  $\mathbf{x}_{max}$  is the row-wise maximum absolute value vector and  $i$ -th entry of  $\mathbf{x}_{max}$  is defined as:  $\{\mathbf{x}_{max}\}_i = \max_{1 \leq j \leq n} |X_{ij}|$ . Proof finished.  $\square$

### A.3 Proof of Proposition 1

PROOF. We assume there come  $p$  edge updates between  $t$  and  $t+1$ . Since we can not access the concrete timestamps of these  $p$  edges, we can assume  $k$ -th each edge arrives at the system at a hidden timestamp  $t + \tau_k$  ( $1 \leq k \leq p$ ) without loss of generality. Specifically, we have  $t + \tau_p = t+1$ . According to [56], the complete evolving process from  $t$  to  $t+1$  can be formulated as:

$$\begin{aligned} M^{(t+1)} &= \bar{\mathcal{A}}M^{(t)} + \bar{\mathcal{B}} \sum_{k=1}^p \sum_{l=0}^{\infty} \alpha \left( (1-\alpha)^l P^{(t+\tau_k)} \right)^l X \Lambda_k \\ &= \bar{\mathcal{A}}M^{(t)} + \bar{\mathcal{B}} \sum_{k=1}^p \alpha \left( I - (1-\alpha)P^{(t+\tau_k)} \right)^{-1} X \Lambda_k. \end{aligned}$$

When setting  $\alpha = \frac{1}{c+1}$ , we have:

$$M^{(t+1)} = \bar{\mathcal{A}}M^{(t)} + \bar{\mathcal{B}} \sum_{k=1}^p \left( I + cL^{(t+\tau_k)} \right)^{-1} X \Lambda_k,$$

where  $L^{(t)}$  is the Laplacian matrix at time  $t$ . Given the exact PPR embedding  $\left( I + cL^{(t+1)} \right)^{-1} X$  at time  $t+1$ , we have:

$$\begin{aligned} &\left\| \left( I + cL^{(t+1)} \right)^{-1} X - \sum_{k=1}^p \left( I + cL^{(t+\tau_k)} \right)^{-1} X \Lambda_k \right\| \\ &= \left\| \left( I + cL^{(t+1)} \right)^{-1} X \sum_{k=1}^p \Lambda_k - \sum_{k=1}^p \left( I + cL^{(t+\tau_k)} \right)^{-1} X \Lambda_k \right\| \\ &= \left\| \sum_{k=1}^p \left( \left( I + cL^{(t+1)} \right)^{-1} - \left( I + cL^{(t+\tau_k)} \right)^{-1} \right) X \Lambda_k \right\| \\ &\leq \left\| \sum_{k=1}^p \left( I + cL^{(t+1)} \right)^{-1} \left( cL^{(t+1)} - cL^{(t+\tau_k)} \right) \left( I + cL^{(t+\tau_k)} \right)^{-1} X \Lambda_k \right\| \\ &\leq c \left\| \sum_{k=1}^p \left( I + cL^{(t+1)} \right)^{-1} \left( L^{(t+1)} - L^{(t+\tau_k)} \right) \left( I + cL^{(t+\tau_k)} \right)^{-1} X \right\| \left\| \sum_{k=1}^p \Lambda_k \right\| \\ &\leq c \sum_{k=1}^p \frac{\left\| \left( L^{(t+1)} - L^{(t+\tau_k)} \right) X \right\|_2}{(1 + c\lambda_1(L^{(t+1)}))(1 + c\lambda_1(L^{(t+\tau_k)}))} \\ &\leq c \sum_{k=1}^p \left\| \left( L^{(t+1)} - L^{(t+\tau_k)} \right) X \right\| \leq \lambda \end{aligned}$$

$\square$

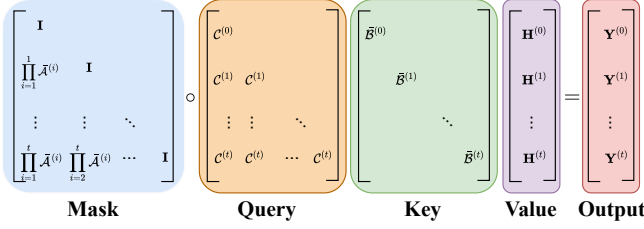
### A.4 Proof of Proposition 2

PROOF. Given  $M^{(0)} = \bar{\mathcal{A}}^{(0)} \cdot \mathbf{0} + \bar{\mathcal{B}}^{(0)} \cdot H^{(0)} = \bar{\mathcal{B}}^{(0)} H^{(0)}$  where  $\bar{\mathcal{A}}^{(t)}$  and  $\bar{\mathcal{B}}^{(t)}$  denote the parameter metrics at time  $t$ . Following the state update principle<sup>3</sup> of Equ. 3, the state matrix at  $t$  can be formulated as:

$$\begin{aligned} M^{(t)} &= \bar{\mathcal{A}}^{(t)} \bar{\mathcal{A}}^{(t-1)} \dots \bar{\mathcal{A}}^{(1)} \bar{\mathcal{B}}^{(0)} H^{(0)} + \bar{\mathcal{A}}^{(t)} \bar{\mathcal{A}}^{(t-1)} \dots \bar{\mathcal{A}}^{(2)} \bar{\mathcal{B}}^{(1)} H^{(1)} + \\ &\quad \dots + \bar{\mathcal{A}}^{(t)} \bar{\mathcal{A}}^{(t-1)} \bar{\mathcal{B}}^{(t-2)} H^{(t-2)} + \bar{\mathcal{A}}^{(t)} \bar{\mathcal{B}}^{(t-1)} H^{(t-1)} + \bar{\mathcal{B}}^{(t)} H^{(t)} \\ &= \sum_{s=0}^t \prod_{i=s+1}^t \bar{\mathcal{A}}^{(i)} \bar{\mathcal{B}}^{(s)} H^{(s)}, \end{aligned} \quad (7)$$

where we define  $\prod_{i=t+1}^t \bar{\mathcal{A}}^{(i)} = I$ . When we vectorize the prediction result over time  $[0, t]$  and restructure the mathematical

<sup>3</sup>Without the loss of generalization, here we simplify the notation and assume the sample time corresponds with the prediction time for a clear presentation.



**Figure 4: The equivalent kernel attention in CODEN.**

expression, we can derive the formula representing our prediction across all time steps:

$$Y^{(0:t)} = (L \odot QK^\top) \cdot V, \quad (8)$$

where (i)  $L$  is a lower-triangular matrix and  $L_{uv} = \prod_{i=v+1}^u \bar{\mathcal{A}}^{(i)}$ ; (ii)  $Q$  is a lower-triangular matrix and  $Q_{uv} = C^{(u)}$ , where  $C$  is the parameters of  $\text{MLP}(\cdot)$ ; (iii)  $K$  is a diagonal matrix and  $K_{uu} = \bar{\mathcal{B}}^{(u)}$ ; (iv)  $V$  is the vectorized node embedding generated at all time steps and  $V(u) = H^{(u)}$ .

As illustrated in Fig. 4, we observe that the fundamental methodology of CODEN for processing temporal information is fully aligned with the definition of the kernel attention mechanism [72], which is developed to mitigate the “over attention” phenomenon [76] in computer vision [74] and natural language processing (NLP) areas [75]. Upon revisiting the motivation behind kernel attention in NLP, where it is crucial to selectively maintain and degrade specific tokens [75], it becomes apparent that the trainable matrices  $\bar{\mathcal{A}}^{(i)}$  ( $0 \leq i \leq t$ ) play a pivotal role.  $\square$

### A.5 Proof of Lemma 2

PROOF. Based on the definition of the Dirichlet Energy and simplifying  $I - A^{(t)\top} D^{(t)-1}$  as  $\Delta$ , we have:

$$\text{DE}(\mathcal{M}_A^{(t)}) = \text{DE} \left( \text{softmax} \left( \frac{(W_q H^{(t)}) \cdot (W_k H^{(s)})^\top}{\sqrt{F}} \right)_{0:t} \cdot (W_v H^{(s)})_{0:t} \right).$$

To simplify the demonstration, we denote

$$S_t = \text{softmax} \left( (W_q H^{(t)}) \cdot (W_k H^{(s)})^\top \right)_{0:t}$$

. Since  $S_t$  is a Row-Stochastic matrix where the sum of each row equals to 1, we can obtain the maximum eigenvalue  $\sigma_{\max}(S_t)$  of  $S_t$  as  $\sigma_{\max}(S_t) \leq 1$ . Then we have:

$$\begin{aligned} \text{DE}(\mathcal{M}_A^{(t)}) &= \text{tr} \left( S_t \left( H^{(s)} \right)_{0:t} \Delta S_t^\top \left( H^{(s)} \right)_{0:t}^\top \right) / F \\ &= \text{tr} \left( \left( H^{(s)} \right)_{0:t} \Delta \left( H^{(s)} \right)_{0:t}^\top S_t S_t^\top \right) / F \\ &\leq \text{tr} \left( \left( H^{(s)} \right)_{0:t} \Delta \left( H^{(s)} \right)_{0:t}^\top \right) \sigma_{\max}(S_t S_t^\top) / F \\ &\leq \text{tr} \left( \left( H^{(s)} \right)_{0:t} \Delta \left( H^{(s)} \right)_{0:t}^\top \right) / F \\ &= \sum_{s=0}^t \text{tr} \left( \left( H^{(s)} \right) \left( H^{(s)} \right)^\top \Delta \right) / F \end{aligned}$$

Following the similar derivation, we can obtain the Dirichlet Energy  $\text{DE}(\mathcal{M}_C^{(t)})$  as:

$$\begin{aligned} \text{DE}(\mathcal{M}_C^{(t)}) &\geq \sum_{s=0}^t \text{tr} \left( \prod_{i=s+1}^t \bar{\mathcal{A}}^{(i)} \bar{\mathcal{B}}^{(s)} \left( H^{(s)} \right) \left( H^{(s)} \right)^\top \Delta \right) \\ &\geq \sum_{s=0}^t \sigma_{\min} \left( \prod_{i=s+1}^t \bar{\mathcal{A}}^{(i)} \bar{\mathcal{B}}^{(s)} \right) \text{tr} \left( \left( H^{(s)} \right) \left( H^{(s)} \right)^\top \Delta \right) \\ &\geq \text{DE}(\mathcal{M}_A^{(t)}), \end{aligned}$$

which holds when  $F \cdot \sigma_{\min}^2 \left( \prod_{i=s+1}^t \bar{\mathcal{A}}^{(i)} \bar{\mathcal{B}}^{(s)} \right) \geq 1$  for  $0 \leq s \leq t$ .  $\square$

### A.6 Proof of Lemma 5

PROOF. According to [32], the estimated vector  $\hat{h}^{(t-1)}$  at time  $t-1$  and the exact PPR vector  $\sum_{l=0}^\infty \alpha \left( (1-\alpha)^l P^{(t-1)} \right)^l \cdot x_i$  have the following relationship:

$$\begin{aligned} \sum_{l=0}^\infty \alpha \left( (1-\alpha)^l P^{(t-1)} \right)^l x_i &= \hat{h}^{(t-1)} + \sum_{l=0}^\infty \alpha \left( (1-\alpha)^l P^{(t-1)} \right)^l r_i \\ &= \hat{h}^{(t-1)} + \alpha \left( I - (1-\alpha)P^{(t-1)} \right)^{-1} r_i \end{aligned}$$

For the purpose of clear demonstration, we express  $(1-\alpha)P^{(t)}$  as  $P^{(t)}$  for the following proof. Then we formulate the output  $\hat{h}^{(t)}$  of Alg. 3 as:

$$\begin{aligned} \hat{h}^{(t)} &= \sum_{j=0}^\infty (P^{(t)})^j (P^{(t)} - P^{(t-1)}) \left( \alpha \sum_{i=0}^\infty (P^{(t-1)})^i x_i - \alpha \sum_{i=0}^\infty (P^{(t-1)})^i r_i^{(t-1)} \right) + \hat{h}^{(t-1)} \\ &= \sum_{j=0}^\infty (P^{(t)})^j (P^{(t)} - P^{(t-1)}) \alpha \sum_{i=0}^\infty (P^{(t-1)})^i x_i - \sum_{j=0}^\infty (P^{(t)})^j (P^{(t)} - P^{(t-1)}) \alpha \sum_{i=0}^\infty (P^{(t-1)})^i r_i^{(t-1)} + \hat{h}^{(t-1)} \end{aligned} \quad (9)$$

The first term of the last line in Equ. 9 can be simplified as:

$$\begin{aligned} &\sum_{j=0}^\infty (P^{(t)})^j (P^{(t)} - P^{(t-1)}) \alpha \sum_{i=0}^\infty (P^{(t-1)})^i x_i \\ &= \sum_{j=1}^\infty (P^{(t)})^j \alpha \sum_{i=0}^\infty (P^{(t-1)})^i x_i - \sum_{j=0}^\infty (P^{(t)})^j \alpha \sum_{i=j}^\infty (P^{(t-1)})^i x_i \\ &= \sum_{j=1}^\infty (P^{(t)})^j \alpha \left( \sum_{i=1}^\infty (P^{(t-1)})^i + I \right) x_i - \left( \sum_{j=1}^\infty (P^{(t)})^j + I \right) \alpha \sum_{i=1}^\infty (P^{(t-1)})^i x_i \\ &= \alpha \sum_{j=1}^\infty (P^{(t)})^j x_i - \alpha \sum_{i=1}^\infty (P^{(t-1)})^i x_i = \alpha \sum_{j=0}^\infty (P^{(t)})^j x_i - \alpha \sum_{i=0}^\infty (P^{(t-1)})^i x_i \end{aligned}$$

Similarly, the second term can be formed as:

$$\begin{aligned} &\sum_{j=0}^\infty (P^{(t)})^j (P^{(t)} - P^{(t-1)}) \alpha \sum_{i=0}^\infty (P^{(t-1)})^i r_i^{(t-1)} \\ &= \alpha \sum_{j=0}^\infty (P^{(t)})^j r_i^{(t-1)} - \alpha \sum_{i=0}^\infty (P^{(t-1)})^i r_i^{(t-1)} \end{aligned}$$

Then Equ. 9 is further expressed as:

$$\begin{aligned}\hat{h}^{(t)} &= \alpha \sum_{j=0}^{\infty} (P^{(t)})^j \mathbf{x}_i - \alpha \sum_{i=0}^{\infty} (P^{(t-1)})^i \mathbf{x}_i - \alpha \sum_{j=0}^{\infty} (P^{(t)})^j \mathbf{r}_i^{(t-1)} + \\ &\quad \alpha \sum_{i=0}^{\infty} (P^{(t-1)})^i \mathbf{r}_i^{(t-1)} + \alpha \sum_{i=0}^{\infty} (P^{(t-1)})^i \mathbf{x}_i - \alpha \sum_{i=0}^{\infty} (P^{(t-1)})^i \mathbf{r}_i^{(t-1)} \\ &= \alpha \sum_{j=0}^{\infty} (P^{(t)})^j \mathbf{x}_i - \alpha \sum_{j=0}^{\infty} (P^{(t)})^j \mathbf{r}_i^{(t-1)}\end{aligned}$$

Since the  $\mathbf{r}_i^{(t-1)}(u) < \epsilon$  for each  $u \in \mathcal{V}^{(t-1)}$ , the proof is finished.  $\square$

### A.7 Proof of Lemma 3

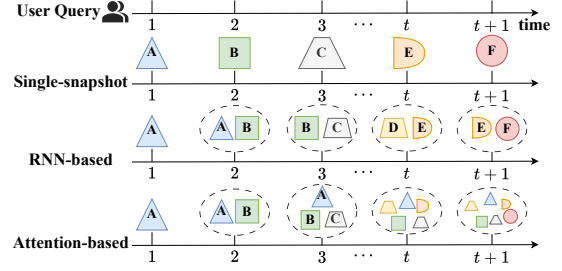
**PROOF.** We first discuss the simple situation where only one edge  $(u, v)$  is added into graph  $\mathcal{G}^{(t)}$  and transfer it into  $\mathcal{G}^{(t+1)}$ . Then we can naturally derive the propagation complexity of a single dimension as  $\|(1 - \alpha)(P^{(t+1)} - P^{(t)})\mathbf{h}^{(t)}\|_1 / \alpha\epsilon$ , where  $\mathbf{h}^{(t)}$  is one dimension of  $H^{(t)}$ . As demonstrated in Alg. 3, given an added edge  $(u, v)$  the residue of  $\|(1 - \alpha)(P^{(t+1)} - P^{(t)})\mathbf{h}^{(t)}\|_1$  can be divided as:

$$\begin{aligned}\|(1 - \alpha)(P^{(t+1)} - P^{(t)})\mathbf{h}^{(t)}\|_1 &= |(1 - \alpha)\mathbf{h}^{(t)}(u) / |\mathcal{N}_{out}^{(t+1)}(u)|| \\ &\quad + \left| \sum_{w \in \mathcal{N}_{out}^{(t+1)}(u)} (1 - \alpha)\mathbf{h}^{(t)}(u) \left( \frac{1}{|\mathcal{N}_{out}^{(t+1)}(u)|} - \frac{1}{|\mathcal{N}_{out}^{(t)}(u)|} \right) \right| \\ &\leq (1 - \alpha)|\mathbf{h}^{(t)}(u)| + \left| |\mathcal{N}_{out}^{(t+1)}(u)| \cdot (1 - \alpha)\mathbf{h}^{(t)}(u) \cdot \frac{1}{|\mathcal{N}_{out}^{(t+1)}(u)|} \right| \\ &\stackrel{(1)}{\leq} 2(1 - \alpha)|\mathbf{h}^{(t)}(u)| + 2(1 - \alpha)\epsilon,\end{aligned}$$

where  $\mathbf{h}^{(t)} = \alpha(I - (1 - \alpha)P^{(t)})^{-1} \mathbf{x}^{(t)}$  and (1) holds since  $\mathbf{h}^{(t)}(u)$  is the underestimation of  $\mathbf{h}^{(t)}(u)$ . Under pattern (i), edges arrive in the system randomly so that each node has the same probability of being the start point of the changed edge. Therefore we can obtain the expected time for each update as  $E[\frac{2(1-\alpha)|\mathbf{h}^{(t)}(u)| + 2(1-\alpha)\epsilon}{\alpha\epsilon}] = \frac{2(1-\alpha)\|\mathbf{x}^{(t)}\|_1}{\epsilon n^{(t)}} + \frac{2(1-\alpha)}{\alpha}$ . Given there exist  $p$  edge updates in a batch at time  $t_k$ , we add the superscript of  $\mathbf{x}^{(t_k)}$  and formulate the complexity for propagation as  $O(p \sum_{i=1}^F \frac{\|\mathbf{x}_i^{(t_k)}\|_1}{\epsilon n^{(t_k)}})$ . Under pattern (ii), we have  $\frac{2(1-\alpha)|\mathbf{h}^{(t)}(u)|}{\alpha\epsilon} \leq \frac{2(1-\alpha) \max_{u \in \mathcal{V}} \|\mathbf{x}^{(t)}(u)\|_1}{\epsilon}$ , which indicates the complexity of this case is  $O(p \sum_{i=1}^F \frac{\max_{u \in \mathcal{V}} \|\mathbf{x}_i^{(t_k)}(u)\|_1}{\epsilon})$   $\square$

### A.8 Proof of Lemma 4

**PROOF.** Based on Alg. 1, for an edge update  $(u, v)$ , the increased mass of  $\sigma$  is at most  $\frac{1-\alpha}{\alpha} \left\| (P^{(t)} - P^{(t-1)}) \cdot \mathbf{x}_{max} \right\|_1 + 2n\epsilon$ . Following the derivation of Lemma 3, we can naturally obtain the upper bound of this term as  $2(1-\alpha)|\mathbf{x}_{max}(u)| + 2n\epsilon$ . Therefore, the largest number of edges contained before  $\sigma > \lambda$  is at least  $\lambda / (2(1-\alpha)|\mathbf{x}_{max}(u)| + 2n\epsilon)$ . Then given  $p$  update events, the number of embeddings sampled  $L$  can be bounded as  $L \leq (2(1-\alpha)|\mathbf{x}_{max}(u)| + 2n\epsilon)p / \lambda$ . Finally, the time complexity of updating the node state in SSM is  $O(n^{(t)}LF^2)$



**Figure 5: The comparison between different TGNN paradigms, where each color denotes a different status of the graph.**

[56, 66]. Therefore, we can naturally obtain the complexity of updating node states as  $O(\frac{\|\mathbf{x}_{max}\|_1 + \epsilon(n^{(t)})^2}{\lambda} pF^2)$  under pattern (i) and  $O(\frac{\max_{u \in \mathcal{V}^{(t)}} \|\mathbf{x}_{max}(u)\|_1 n^{(t)} + \epsilon(n^{(t)})^2}{\lambda} pF^2)$  under pattern (ii).  $\square$

### A.9 Related Work

**A.9.1 Single-snapshot Methods.** The single-snapshot GNN methods directly utilize the node embedding as the node states, which will be rapidly updated based on the new interactions. Specifically, single-snapshot GNN methods aim to transmit the updates into node embeddings  $H^{(t+1)}$  from  $H^{(t)}$  with the minimal time cost following:

$$H^{(t+1)} = \text{MSG}(\mathcal{G}^{(t)}, H^{(t)}, e_{t+1}). \quad (10)$$

To facilitate the incremental updating based on these pre-computed node embeddings, Instant [32], DynAnom [33] and IDOL [41] explore the invariant rules of graph propagation and conduct an invariant-based algorithm of Personalized PageRank (PPR) to refresh the node embedding locally. The dominant update complexity of these methods can be bounded as  $O(p \sum_{i=1}^F \frac{\|\mathbf{x}_i^{(t)}\|_1}{\epsilon n^{(t)}})$  given  $p$  updates and  $\epsilon$  approximation error [32]. However, only focusing on the current snapshot will miss significant interactions of past time steps, leading to sub-optimal prediction results.

**A.9.2 RNN-based Methods.** RNN-based methods are generally based on the classical RNN [52], GRU [53], LSTM [54] algorithm, etc., which simply use both the current input and the previous hidden state to iteratively capture temporal dependencies. For example, the GRU algorithm updates the node states at time  $t + 1$  as:

$$\begin{aligned}Z^{(t+1)} &= \text{sigmoid}(W_Z H^{(t+1)} + U_Z M^{(t)} + B_Z) \\ R_t &= \text{sigmoid}(W_R H^{(t+1)} + U_R M^{(t)} + B_R) \\ \tilde{H}^{(t+1)} &= \tanh(W_H H^{(t+1)} + U_H (R_t \circ M^{(t)})) + B_H \\ M^{(t+1)} &= (1 - Z^{(t+1)}) \circ M^{(t)} + Z^{(t+1)} \circ \tilde{H}^{(t+1)},\end{aligned}$$

where  $W, U, B$  denote the trainable parameters of the linear layer and  $H^{(t+1)}$  can be updated using the  $\text{MSG}(\cdot)$ . Specifically, TGCN [43] and EvolveGCN [34] and incorporate the Graph Convolutional Network (GCN)[55] as the  $\text{MSG}(\cdot)$  function to regenerate the node embeddings while coupling with an RNN-based module to learn temporal node representations. Similarly, MPNN [35] transforms the node embeddings at different time steps into an RNN module

and then captures the long-range dependency in the final hidden state. This category of method generally requires  $O(Km^{(t)}F)$  and  $O(pn^{(t)}F^2)$  to update the node embeddings and states <sup>4</sup>, respectively. Due to their iterative structure, RNN-based methods can efficiently update node states. However, the simplistic recurrence mechanism often leads to difficulties in retaining historical information, especially as the graph size and temporal scope expand [56, 57].

**A.9.3 Attention-based Methods.** To address the forgetting problem of RNNs, attention-based methods rely on the attention mechanism and abstain from using recurrence form, which encodes the position of sequences and enables the efficient information flow from past to current representations. We take the representative work APAN [83] as an example to demonstrate the core mechanism of this category. Considering matrices  $Q \in \mathbb{R}^{n \times F}$  denoted as "query",  $K \in \mathbb{R}^{n \times F}$  denoted as "keys", and  $V \in \mathbb{R}^{n \times F}$  denoted as "values", the classical attention algorithms perform the following computation to obtain the optimized embeddings:

$$M^{(t+1)} = \text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{F}}\right)V,$$

$$Q = H^{(t+1)}W_q, K = M^{(t)}W_k, V = M^{(t)}W_v,$$

where  $W_q, W_k, W_v \in \mathbb{R}^{F \times F'}$  are the network parameters. The dot-product term  $\left(\frac{QK^T}{\sqrt{F}}\right)$  takes the role of weighting the interactions between entity "query-key" pairs. A higher value within this term increases the contribution of  $V$  to the embedding space. Thus, attention-based methods create the expressive attention score to capture the relationship between the current embedding and the state of the last time step. Following this intuition, DySat [37] employs the generalized GAT module [58] to integrate the embeddings of a single node from different time steps to generate its refreshed one. ASTGCN [38] and DNNTSP [45] further incorporate the attention mechanism to capture the spatial and temporal dependency for enhanced representation quality. For each time step, these methods generally require  $O(T \left(n^{(t)}\right)^2 F)$  time complexity to calculate the final output given  $T$  time step. While attention mechanisms can retain the most relevant parts of the sequences to avoid the forgetting issue, they can become computationally expensive with frequent updates [59].

**A.9.4 Other TGNN methods.** (i) *SNN-based methods.* Another typical mechanism of this category is based on the biological Spiking Neural Networks (SNNs), which simulate the brain behaviors and maintain the membrane potential given the data sequences. SpikeNet [1] retrieves the node embeddings from multiple time steps and finally generates the prediction results by the spike firing process. Dy-SIGN [84] incorporates SNNs mechanism into dynamic graphs to mitigate the information loss and memory consumption problem. Nevertheless, the demand for multiple simulation steps to generate reliable embeddings can significantly degrade the efficiency of these TGNN methods. (ii) *SSM-based methods.* There are also some works which employs the SSM mechanism on temporal graphs. For example, STG-Mamba [69] formulates the feature stream of each node as the long-term context, which improves the

embedding quality for feature-varying graphs. Graph-SSM [56] addresses the unobserved graph mutations between consecutive snapshots, and achieves an effective discretization with long-term information. However, these methods only focus on the discrete-time dynamic graph and fail to model the continuous topology changes as the graph evolves. Furthermore, directly adapting these methods will incur significant computational overhead as the graph evolves, creating a gap between current algorithms and continuous prediction in practical dynamic scenarios.

**A.9.5 State Space Models.** In recent years, structured State Space Models (S4) [68, 85] have emerged as promising frameworks for modeling long-distance sequences, offering the advantage of only a linear increase in computational cost. Given the input series  $x(t) \in \mathbb{R}$ , structured State Space Models (S4) [68, 85] formulates the state variable series  $h(t) \in \mathbb{R}$  and the output series  $y(t) \in \mathbb{R}$  using the following equations:

$$\begin{cases} h'(t) = \mathcal{A}h(t) + \mathcal{B}x(t) \\ y(t) = \mathcal{C}h(t), \end{cases} \quad \begin{cases} h_t = \tilde{\mathcal{A}}h_{t-1} + \tilde{\mathcal{B}}x_t \\ y_t = \mathcal{C}h_t, \end{cases} \quad (11)$$

where  $(\mathcal{A}, \mathcal{B}, \mathcal{C}, \tilde{\mathcal{A}}, \tilde{\mathcal{B}})$  are trainable parameters. Here Equ. 11 and 12 are the continuous and discretization process, respectively. Compared with the RNN-based or attention-based mechanism, S4's strength lies in its adherence to a linear mechanism, which guarantees enhanced stability control [56]. This facilitates effective long-term modeling of sequences through meticulous initialization of state space layer parameters [86, 87]. As a result, we employ the SSMs as our temporal-processing unit to meticulously balance the accuracy and efficiency of our CODEN framework.

**A.9.6 Forward Push.** The prevailing paradigm for existing PPR-based propagation is derived from the concept of *Forward Push*, which calculates the approximated solution  $z_i = \sum_{l=0}^{\infty} \alpha(1-\alpha)^l \left(A^{(t)}{}^T D^{(t)-1}\right)^l x_i^{(t)}$  ( $1 \leq i \leq F$ ) under a given error bound  $\epsilon$ , where  $\alpha$  is the decay factor of random walk. Specifically, *Forward Push* (depicted in Alg. 2) assigns the attribute vector  $x_i^{(t)}$  to the *residue vector*  $r^{(t)}$  (e.g.,  $r^{(t)} = x_i^{(t)}$ ), which represents the unpropagated mass of attribute vector  $x_i^{(t)}$ .<sup>5</sup> We iteratively conduct the following two steps: (1) For each node  $s \in \mathcal{V}^{(t)}$  such that  $r^{(t)}(s) > \epsilon$ ,  $(1-\alpha)$  fraction of  $r^{(t)}(s)$  will be propagated into the out-neighbors  $t \in \mathcal{N}_{out}(s)$  averagely. (2)  $(1-\alpha)$  fraction of  $r^{(t)}(s)$  will be transferred into the reserve vector  $h^{(t)}$  and then  $r^{(t)}(s)$  is set as 0. This iteration will be terminated until  $r^{(t)}(s) \leq \epsilon$  for all nodes  $s \in \mathcal{V}^{(t)}$  and we can deploy  $h^{(t)}$  as the approximated node embedding which satisfies  $|h^{(t)}(s) - z^{(t)}(s)| \leq \epsilon$  for each node  $s \in \mathcal{V}^{(t)}$ .

## A.10 Introduction of baseline methods

**A.10.1 Single-snapshot methods.** **Instant** [32], **DynAnom** [33], **IDOL** [41]. These three methods inherit the updating skeleton of PPR [88] to incrementally update the node embeddings. Following the analysis of [32], the dominant complexity of updating embeddings can be formulated as  $O\left(\frac{\|x^{(t)}\|_1}{\epsilon n}\right)$  for a feature vector  $x^{(t)}$

<sup>4</sup>For a clear presentation, we assume the dimension of node state  $F' = F$

<sup>5</sup>In the following sections, we will omit the subscript  $i$  for simplicity.

---

**Algorithm 2: Forward Push**

---

**Input** : Graph  $\mathcal{G}^{(t)} = (\mathcal{V}^{(t)}, \mathcal{E}^{(t)})$ , reserve vector  $\mathbf{h}^{(t)}$ ,  
residue vector  $\mathbf{r}^{(t)}$ .  
**Output**: Reserve vector  $\mathbf{h}^{(t)}$   
1  $\mathbf{h}^{(t)} = \mathbf{h}^{(t-1)}$ ;  
2 **while** exists  $s \in \mathcal{V}$  such that  $\mathbf{r}^{(t)}(s) > \epsilon$  **do**  
3     **foreach**  $v \in \mathcal{N}_{out}^{(t)}(s)$  **do**  
4          $\mathbf{r}^{(t)}(v) += (1 - \alpha) \cdot \frac{\mathbf{r}^{(t)}(s)}{|\mathcal{N}_{out}^{(t)}(s)|}$ ;  
5      $\mathbf{h}^{(t)}(s) += \alpha \cdot \mathbf{r}^{(t)}(s)$ ;  $\mathbf{r}^{(t)}(s) = 0$ ;

---

assuming  $\alpha$  as the constant. Hence we summarize the complexity of this category as  $O\left(p \sum_{i=1}^F \|\mathbf{x}_i^{(t)}\|_1\right)$  for  $F$  dimensions given  $\epsilon = \frac{1}{n}$ .

**A.10.2 RNN-based methods.** TGN [42]. TGN proposes a general framework for learning the representation in CTGD, which models the past interactions between nodes using a compressed node state vector. Based on existing techniques, TGN provides flexible modules to compute the embeddings and update node states. For example, TGN utilizes the attention mechanism and GCN [55] to compute the node embeddings and LSTM or GRU module to update the node state.

TCGN [43], EvolveGCN [34], MPNN [35], ROLAND [44]. These four methods employ a unified pipeline to update the node embeddings and states. Specifically, these methods conduct the graph propagation based on GCN [55] and update the state with GRU or LSTM units. Since each update will require to compute the embeddings from scratch, we summarize their complexity as  $O(Km^{(t)}F)$  and  $O(pn^{(t)}F^2)$  given  $p$  updates.

**A.10.3 Attention-based methods.** DySat [37], ASTGCN [38]. DySat and ASTGCN utilize structural and temporal attention mechanism for dynamic graph representation learning. Both of them employ graph propagation by GCN. As a result, the graph propagation and state update require  $O(n^2F)$  time complexity for each update.

TGAT [24]. TGAT adopts the GraphSage [60] for the embedding computation. Different from DySat, TGAT aims to formulate the interactions between the time encoding and node features. However, due to the polynomial property of the attention mechanism, TGAT still needs to consume  $O(n^2F)$  time for state update.

DNNTSP [45], SEIGN[46], DyGFormer [47]. These three methods follow a similar intuition to our alternative model CODEN-A, which intends to calculate the interactions between different time steps for each node. Given the time step  $t$ , this framework needs to consume  $O(tn^2F)$  time for the state update.

### A.11 Accuracy-guaranteed Embedding Update

As we mentioned, our objective is to efficiently propagate frequent updates into the node embeddings from the preceding time step. The rationale behind this objective is that not all node embeddings can be significantly influenced by the updates, since an edge update causes only a minimal impact on distant nodes. However, current state-of-the-art algorithms [1, 42, 49] intricately intertwine graph convolution with node states utilizing trainable parameters in the message function  $\text{MSG}(\cdot)$  of Equ. 1, requiring either repetitive convolution operations or a complete recalculation of states upon updates. To address these issues, we propose to incrementally

---

**Algorithm 3: Embedding Update**

---

**Input** : Graph  $\mathcal{G}^{(t)} = (\mathcal{V}^{(t)}, \mathcal{E}^{(t)})$ , update events  
 $\Gamma = \{(u_1, v_1), (u_2, v_2), \dots, (u_p, v_p)\}$ , reserve vector  
 $\mathbf{h}^{(t)}$ , old adjacent matrix  $\mathbf{P}^{(t)}$ , updated adjacent matrix  
 $\mathbf{P}^{(t+p)}$ .  
**Output**: Reserve vector  $\mathbf{h}^{(t+p)}$   
1  $\mathbf{r}^{(t+p)} = \mathbf{0}$ ;  
2  $\mathcal{V}_{changed} \leftarrow \{u \mid \text{the nodes whose out degree has changed}\}$ ;  
3 **foreach**  $u \in \mathcal{V}_{changed}$  **do**  
4      $\mathcal{N}_{add}(u) \leftarrow \{v \mid \text{the added neighbors of } u\}$ ;  
5      $\mathcal{N}_{del}(u) \leftarrow \{v \mid \text{the deleted neighbors of } u\}$ ;  
6     **foreach**  $v \in \mathcal{N}_{add}(u)$  **do**  
7          $\mathbf{r}^{(t+p)}(v) += (1 - \alpha)\mathbf{h}^{(t)}(u) / |\mathcal{N}_{out}^{(t+p)}(u)|$ ;  
8     **foreach**  $v \in \mathcal{N}_{del}(u)$  **do**  
9          $\mathbf{r}^{(t+p)}(v) -= (1 - \alpha)\mathbf{h}^{(t)}(u) / |\mathcal{N}_{out}^{(t)}(u)|$ ;  
10    **foreach**  $w \in \mathcal{N}_{out}^{(t+p)}(u) \setminus (\mathcal{N}_{add}(u) \cup \mathcal{N}_{del}(u))$  **do**  
11          $\mathbf{r}^{(t+p)}(w) += (1 - \alpha)\mathbf{h}^{(t)}(u) \left( \frac{1}{|\mathcal{N}_{out}^{(t+p)}(u)|} - \frac{1}{|\mathcal{N}_{out}^{(t)}(u)|} \right)$ ;  
12 **Forward Push** ( $\mathcal{G}^{(t+p)}, \mathbf{h}^{(t)}, \mathbf{r}^{(t+p)}$ );

---

integrate update events into the node embeddings  $\mathbf{h}^{(t)}$  in the subsequent sections. This step will then trigger an update to the node states  $\mathbf{M}^{(t)}(u)$  to effectively document the ongoing evolutionary process.

**Compensated Propagation.** In this section, we discuss how to efficiently update the node embedding  $\mathbf{H}^{(t)}$  and let  $\|\mathbf{H}^{(t)} - \mathbf{Z}^{(t)}\|_1 \leq n^{(t)}\epsilon$  holds for each time step  $t$  as demonstrated in 1. Despite of the recent advancements to incrementally update the PPR-based embeddings [32, 33], it is pointed out that this invariant-based adjustment can not provide the accuracy bound after updating [65, 89]. To address this challenge, we draw on insights from recent works on Random Walk with Restart (RWR) [65] to compensate for embedding distance across different time steps.

Without loss of generality, we start with an example that  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} \in \mathbf{X}^{(t)}$  and a new edge  $e_{t+1} = (u, v)$  is added at time  $t + 1$ , creating an unbounded value difference between the old embedding  $\mathbf{h}^{(t)}$  and the approximation target  $\mathbf{z}^{(t+1)}$ . Note that  $\mathbf{h}^{(t)}$  is the approximation of  $\mathbf{z}^{(t)}$ , and we have:

$$\begin{aligned} \mathbf{z}^{(t+1)} - \mathbf{z}^{(t)} &= \alpha \left( \left( \mathbf{I} - (1 - \alpha)\mathbf{P}^{(t+1)} \right)^{-1} - \left( \mathbf{I} - (1 - \alpha)\mathbf{P}^{(t)} \right)^{-1} \right) \mathbf{x}^{(t+1)} \\ &\stackrel{(1)}{=} \alpha(1 - \alpha) \left( \mathbf{I} - (1 - \alpha)\mathbf{P}^{(t+1)} \right)^{-1} \Delta \mathbf{P} \left( \mathbf{I} - (1 - \alpha)\mathbf{P}^{(t)} \right)^{-1} \mathbf{x}^{(t+1)} \\ &= \underbrace{\alpha \left( \mathbf{I} - (1 - \alpha)\mathbf{P}^{(t+1)} \right)^{-1}}_{\text{propagation process}} \cdot \underbrace{(1 - \alpha)\Delta \mathbf{P} \mathbf{z}^{(t)}}_{\text{compensated vector}} \end{aligned} \quad (13)$$

where (1) holds since  $\mathbf{U}^{-1} - \mathbf{V}^{-1} = \mathbf{U}^{-1}(\mathbf{V} - \mathbf{U})\mathbf{V}^{-1}$  and  $\Delta \mathbf{P} = (\mathbf{P}^{(t+1)} - \mathbf{P}^{(t)})$ . Based on Equ. 13, we have a key observation: *the difference between node embedding  $\mathbf{z}^{(t)}$  and  $\mathbf{z}^{(t+1)}$  can be compensated by propagating the new feature vector  $(1 - \alpha)(\mathbf{P}^{(t+1)} - \mathbf{P}^{(t)})\mathbf{z}^{(t)}$  along the updated graph  $\mathcal{G}^{(t+1)}$ .* Since  $\mathbf{h}^{(t)}$  is the approximation of  $\mathbf{z}^{(t)}$  with the allowed error, we hence propose to implement the embedding update based on the *compensated vector*  $(1 - \alpha)\Delta \mathbf{P} \mathbf{h}^{(t)}$ .

**Scalable Batch Update.** By adhering to the principle above, we extend the incremental update of node embedding into the batch

setting, which is depicted in Alg. 3. Consider the graph  $\mathcal{G}^{(t)}$  at time  $t$  and the batch update events in  $\Gamma = \{(u_1, v_1), (u_2, v_2), \dots, (u_p, v_p)\}$  containing  $p$  edge updates, where  $(u_i, v_i)$  ( $1 \leq i \leq p$ ) will be viewed as deletion if it exists in  $\mathcal{G}^{(t+i-1)}$  otherwise as an addition. First, we obtain the set  $\mathcal{V}_{\text{changed}}$  denoting the nodes whose out-degree has changed. Given a node  $u \in \mathcal{V}_{\text{changed}}$  and for each node  $v \in \mathcal{N}_{\text{out}}^{(t)}(u) \cup \mathcal{N}_{\text{out}}^{(t+p)}(u)$ , the corresponding entry of  $\Delta P = P^{(t+p)} - P^{(t)}$  can be expressed as:

$$\Delta P[v, u] = \begin{cases} \frac{1}{|\mathcal{N}_{\text{out}}^{(t+p)}(u)|}, & v \text{ is the added neighbor of } u, \\ -\frac{1}{|\mathcal{N}_{\text{out}}^{(t)}(u)|}, & v \text{ is the deleted neighbor of } u, \\ \frac{1}{|\mathcal{N}_{\text{out}}^{(t+p)}(u)|} - \frac{1}{|\mathcal{N}_{\text{out}}^{(t)}(u)|}, & \text{otherwise.} \end{cases}$$

Then, we multiply  $\Delta P[v, u]$  with  $(1-\alpha)\mathbf{h}^{(t)}(u)$  and finish the computation of the compensated vector  $(1-\alpha)(P^{(t+p)} - P^{(t)})\mathbf{h}^{(t)}$  (lines 6-11). Moreover, we assign this new feature vector as the unpropagated residue vector  $\mathbf{r}^{(t+p)}$ . Note that  $\mathbf{r}^{(t+p)}$  may exceed the permissible error for certain nodes, such as when  $\mathbf{r}^{(t+p)}(u) > \epsilon$  ( $u \in \mathcal{V}^{(t+p)}$ ). Consequently, we trigger the *Forward Push* mechanism [79] again to propagate the residue vector  $\mathbf{r}^{(t+p)}$  (line 12), which will subsequently influence the embedding values of other nodes and ensure the desired accuracy. Additionally, our algorithm can also be compatible with the attribute changes, which is detailed in the Appendix [40].

**Theoretical Accuracy Guarantee.** Alg. 3 returns the updated embeddings by propagating the compensated vector along the updated graph, where the compensated vector is derived using  $\mathbf{z}^{(t)}$ , as specified in Equ. 13. Although we process this vector with the approximated vector  $\mathbf{h}^{(t)}$ , we can still establish an error bound for the output vector  $\mathbf{h}^{(t+p)}$ , which is the result of propagation from scratch at time  $t+p$ . We formally state this accuracy guarantee in the following lemma:

**LEMMA 5.** *Given the normalized adjacent matrix  $P^{(t)}$  at time  $t$  and the update event set  $\Gamma = \{(u_1, v_1), (u_2, v_2), \dots, (u_p, v_p)\}$ , Alg 3 can output the approximated embedding vector  $\mathbf{h}^{(t+p)}$  which satisfies:*

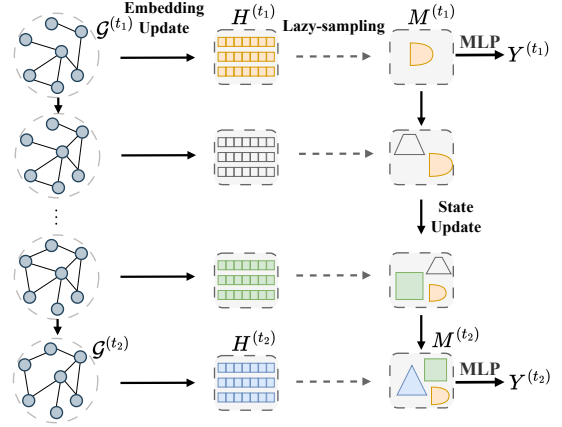
$$\|\mathbf{H}^{(t+p)} - \mathbf{Z}^{(t+p)}\|_1 \leq n^{(t)}\epsilon$$

where  $\mathbf{Z}^{(t+p)} = \sum_{l=0}^{\infty} \alpha(1-\alpha)^l (P^{(t+p)})^l \mathbf{X}^{(t+p)}$ .

**Comparison with the invariant-based scheme.** Compared with the invariant-based scheme in existing related works [32, 33, 41], we claim that our method stands out as the first to offer guaranteed accuracy of updated embeddings in dynamic settings. The rationale for this difference is that our method offers a more stable adjustment by leveraging a compensate vector, making it more robust in scenarios with significant node degree changes. A running example underpinning this insight is further provided in the Appendix [40]. In summary, our accuracy-guaranteed embedding update forms the novel theoretical foundation of CODEN (e.g., Lemma 1 and Proposition 1), thereby conferring substantial benefits in producing high-quality representations.

## A.12 Extension to attribute changes.

Note that we have made an assumption that  $\mathbf{x}^{(t+p)} = \mathbf{x}^{(t)}$  such that we can focus on the topological change of the graph. Actually,



**Figure 6: An illustration of the CODEN framework.**  $M^{(t)}$ ,  $H^{(t)}$  and  $Y^{(t)}$  denote the node state matrix, embedding matrix and the prediction results at time  $t$ , respectively.

this assumption can be easily relaxed since extending our approach to dynamic-attribute graphs is much more straightforward. Specifically, given  $\Delta \mathbf{x} = \mathbf{x}^{(t+p)} - \mathbf{x}^{(t)}$  and  $P^{(t+p)} = P^{(t)}$ , we can naturally derive the difference between  $\mathbf{z}^{(t)}$  and  $\mathbf{z}^{(t+p)}$  as:

$$\mathbf{z}^{(t+p)} - \mathbf{z}^{(t)} = \alpha \left( I - (1-\alpha)P^{(t+1)} \right)^{-1} \Delta \mathbf{x}. \quad (14)$$

Therefore, we can view  $\Delta \mathbf{x}$  as the compensated vector and easily obtain the updated embedding by directly invoking *Forward Push* ( $\mathcal{G}^{(t+p)}, \mathbf{h}^{(t)}, \Delta \mathbf{x}$ ).<sup>6</sup>

## A.13 Discussion of the Invariant-based Scheme

Considering  $p$  edge updates, existing related works [32, 33, 41] adhere the invariant relationship in PPR and perform the following rules to locally update the node embedding:  $\mathbf{r}^{(t+p)}(u) \leftarrow \mathbf{r}^{(t)}(u) - \frac{\Delta_p(u)\mathbf{h}^{(t)}(u)}{\alpha|\mathcal{N}_{\text{out}}^{(t)}(u)|}$ ,  $\mathbf{r}^{(t+p)}(v) \leftarrow \mathbf{r}^{(t)}(v) + \frac{(1-\alpha)\Delta_p(u)\mathbf{h}^{(t)}(u)}{\alpha|\mathcal{N}_{\text{out}}^{(t)}(u)|}$ ,  $\mathbf{h}^{(t+p)}(u) \leftarrow \mathbf{h}^{(t)}(u) \cdot \frac{|\mathcal{N}_{\text{out}}^{(t+p)}(u)|}{|\mathcal{N}_{\text{out}}^{(t)}(u)|}$ , where  $\Delta_p(u) = |\mathcal{N}_{\text{out}}^{(t+p)}(u)| - |\mathcal{N}_{\text{out}}^{(t)}(u)|$ . The *Forward Push* process is then repeated to ensure that local adjustments propagate to distant nodes. Although these methods follow a similar pipeline to CODEN, we emphasize that our initialization process prior to *Forward Push* is fundamentally different.

To illustrate this distinction, we present a toy example in Fig. 7 where 3 edges are added between node  $v_0$  and  $v_2 \sim v_4$ . Based on this result, we see that the invariant-based scheme can not guarantee the accuracy of node embeddings since  $|\mathbf{h}(v) - \mathbf{z}(v)| = |0.75 - 0.08| > 0.5 = \epsilon$  for  $v = v_2 \sim v_4$ . In contrast, our accuracy-based can still maintain the embedding accuracy after the *Forward Push* algorithm, which confers substantial benefits in producing high-quality representations for later node state formulation.

## A.14 Additional Experimental Results

**A.14.1 Experimental Setting.** Since most of the compared RNN-based and attention-based methods are designed for small-scale graphs, we adopt the neighboring sampling techniques [60] for these methods when applied on the *Reddit* and *Patent* datasets to

<sup>6</sup>In the remaining sections, we mainly focus on the topology change and assume  $\mathbf{X}^{(t)}$  remains unchanged in the following sections for a clear presentation.



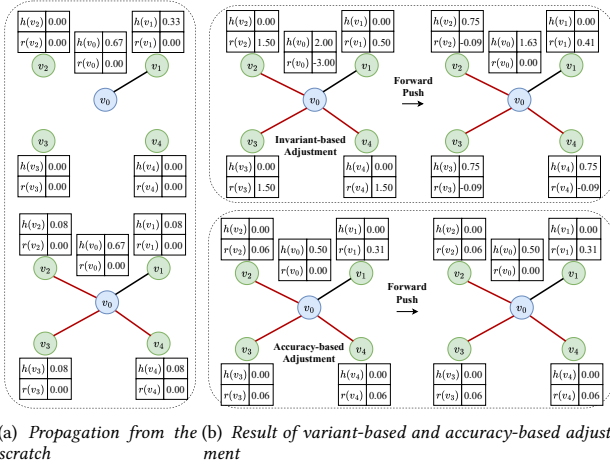


Figure 7: A toy example to compare the results of 4 kinds of push operations when we set  $\alpha = 0.5$  and  $\epsilon = 0.5$ . (a) shows the initial push result before the update. (b) shows the ground truth result of the push from scratch. (c) shows the result of the variant-based push, where the error of the embeddings of node  $v_2 \sim v_4$  have exceeded the bound  $\epsilon$ . (d) shows the result of the accuracy-based push, where the accuracy of node embeddings is guaranteed.

Table 6: Parameter settings. Here "lr" means the learning rate, "K" means the number of convolution layers, "hidden" means the hidden size of the network, and "batch number" means the number of neighbor sampling for baselines.

Datasets	lr	K	hidden	batch number
DBLP	1e-3	4	1024	12
Tmall	1e-3	4	1024	12
Reddit	1e-3	4	512	12
Patent	1e-3	4	512	12

avoid the out-of-memory problem. We summarize the experimental settings for all baselines in Tab. 6, where the common parameters, such as learning rate and hidden size, are maintained consistently for CODEN.

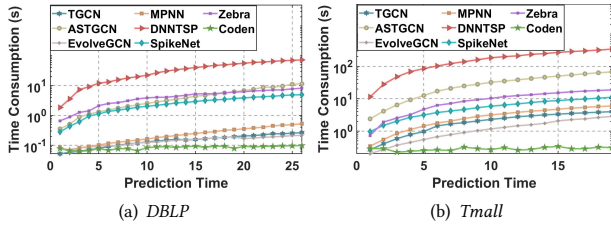


Figure 8: Average training time (per epoch) for each prediction time.

#### A.14.2 Training Time on DBLP and Tmall.

A.14.3 Information compression under flexible arrival patterns. In the experiments above, we added edges to the initial graph following an edge-inclined arrival pattern. However, using a more flexible

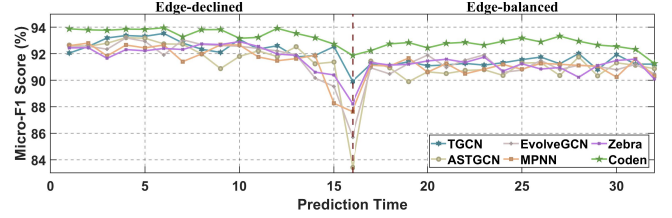


Figure 9: Micro-F1 scores under the Edge-declined and Edge-balanced patterns on *Reddit* dataset.

arrival pattern would allow for a more thorough examination of the models' memory capabilities, particularly in assessing whether they experience significant performance degradation. Therefore, we employ more edge arrival patterns on each model and report the corresponding accuracy on *Reddit* dataset, as depicted in Fig. 9. In the edge-declined pattern, we reverse the order of the edge-inclined setting and remove the batch of edges sequentially from the final snapshot. Then, we randomly add and remove one batch of edges simultaneously to the initial graph, which is denoted as the edge-balanced pattern. In the edge-declined pattern, we observe that as the edges are missing gradually, all methods suffer from a performance decrease. However, CODEN is only slightly affected by the removal of edge batches, which particularly achieves a non-trivial improvement on the initial graph compared with the edge-inclined pattern at time step 16 (e.g., 91.86% vs 87.86%). Moreover, CODEN demonstrates a significant improvement in this pattern, achieving an average prediction accuracy of 93.48%, indicating its ability to effectively compress historical information for future predictions.

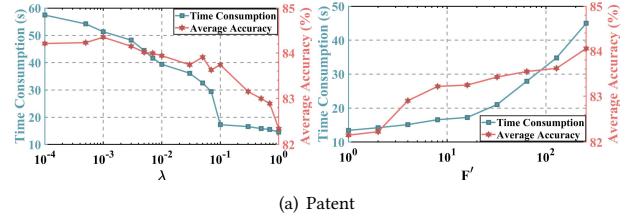


Figure 10: The average training time and the average accuracy when setting different  $\lambda$  and  $F'$  on *Patent* dataset.